

CENTER FOR
RIBBON
CALCULATIONS
LYSIS CBS

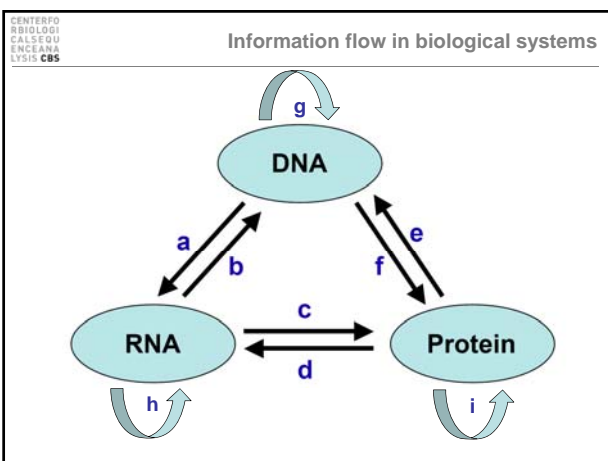
DNA as Biological Information

Rasmus Wernersson
Henrik Nielsen

CENTER FOR
RIBBON
CALCULATIONS
LYSIS CBS

Overview

- Learning objectives
 - Biological information
 - DNA sequencing techniques and DNA data
 - File formats used for biological data
 - Introduction to the GenBank database



[illegible]

PCR

PCR

PCR

Denaturation
96°C, 30 sec

Annealing
~55°C, 30 sec

Extension
72°C, 30 sec

35 cycles

DNA Sequence:
5'-CTGAGTATGAGACCTATAGGTACGGTGGCCATTCTGTCTGATCCCGGACTACTACAGAA-3'
3'-GAGCATGATGG-5'

CENTERFO
BRILOGI
CALSEQU
ENCEAN
LYSIS **CBS**

PCR

AMOUNT OF DNA

PCR CYCLE NUMBER

Real target

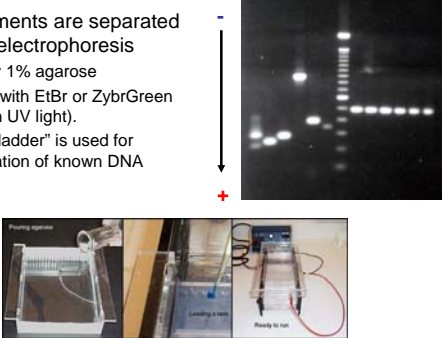
Single-primer target (500)

Single-primer target (1000)

Animation: <http://www.people.virginia.edu/~rjh9u/pcranim.html>
PCR graph: <http://pathmicro.med.ac.edu/ncr/realtime-home.htm>

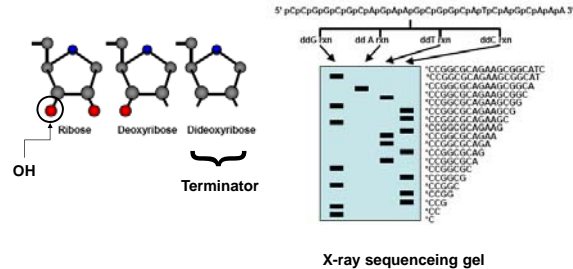
Gel electrophoresis

- DNA fragments are separated using gel electrophoresis
 - Typically 1% agarose
 - Colored with EtBr or ZylbrGreen (glows in UV light).
 - A DNA "ladder" is used for identification of known DNA lengths.



Gel picture: <http://www.pharmaceutical-technology.com/projects/roche/images/roche3.jpg>
 PCR setup: <http://arbl.cvmbs.colostate.edu/bbooks/genetics/biotech/gels/agardna.html>

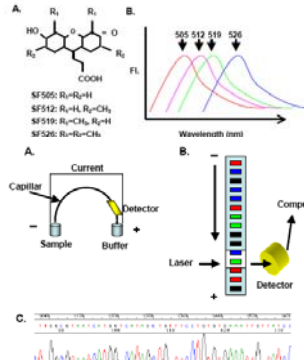
The Sanger method of DNA sequencing



Images: http://www.idtdna.com/support/technical/TechnicalBulletinPDF/DNA_Sequencing.pdf

Automated sequencing

- The major break-through of sequencing has happened through *automation*.
- Fluorescent dyes.
- Laser based scanning.
- Capillary electrophoresis
- Computer based base-calling and assembly.

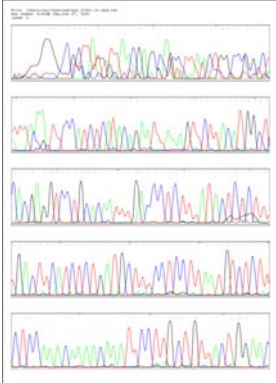


Images: http://www.idtdna.com/support/technical/TechnicalBulletinPDF/DNA_Sequencing.pdf

CENTER FOR
BIOLOGICAL
CALSEQU
ENGINEERING
LYSIS CBS

Handout exercise: "base-calling"

- Handout: Chromatogram
- Groups of 2-3.
- Tasks:
 - Identify "difficult" regions
 - Identify likely errors
 - Try to estimate the best interval to use



CENTER FOR
BIOLOGICAL
CALSEQU
ENGINEERING
LYSIS CBS

Sequence read mapping



CENTER FOR
BIOLOGICAL
CALSEQU
ENGINEERING
LYSIS CBS

DNA sequencing – history

1972 Recombinant DNA technology [Paul Berg].

1976 The first sequenced genome, the bacteriophage MS2 [Walter Fiers *et al.*]

1977 DNA sequencing by chemical cleavage [Allan Maxam & Walter Gilbert]; DNA sequencing by enzymatic synthesis [Fred Sanger].

1982 *GenBank* (public data base of DNA sequences).

1987 The first automatic sequencing machine, *Prism 373* [Applied Biosystems].

1990 The *Human Genome Project* is launched.

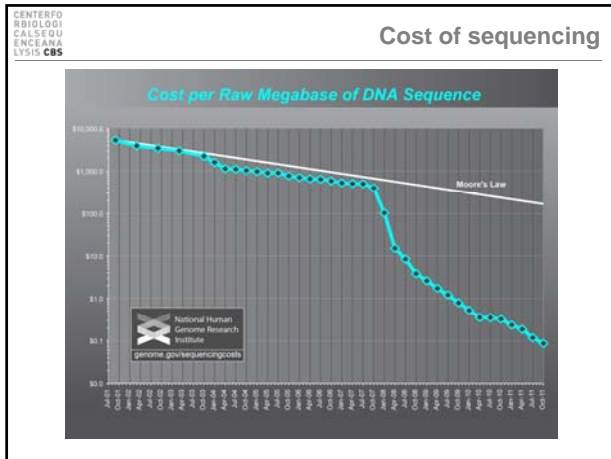
1995 The first genome of a free-living organism, the bacterium *Haemophilus influenzae* (1.8 Mb) [The Institute for Genomic Research (TIGR)].

1996 The first genome of a eukaryote, Baker's yeast, *Saccharomyces cerevisiae* (12.1 Mb) [International consortium].

1998 The first genome of an animal, the worm *Caenorhabditis elegans* (97Mb) [Sanger Center and partners].

2001 The first "drafts" of the human genome (3Gb) [Human Genome Project Consortium (Nature, 15 Feb) + Celera (Science, 16 Feb)].

April 11, 2013 *GenBank release 195* contains 164,136,731 sequences with a total of 151,178,979,155 nucleotides (the files take up 594 GB).



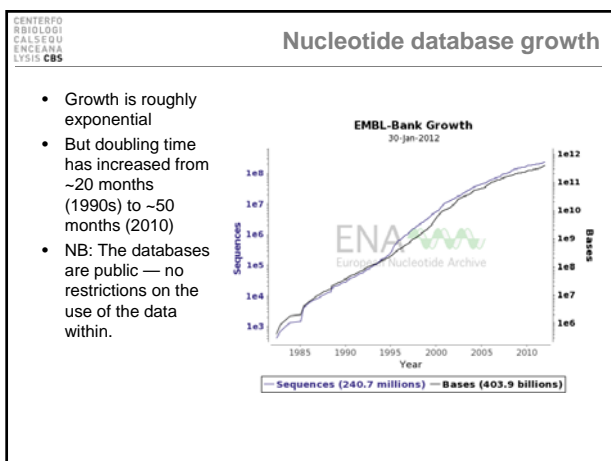
CENTER FOR BIOLOGICAL CHEMISTRY
ENGLAND
LYSIS CBS

Background - Nucleotide databases

- **GenBank**, <http://www.ncbi.nlm.nih.gov/Genbank/>
- National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National Institutes of Health (NIH), USA
- Established in 1982.
- **EMBL**, <http://www.ebi.ac.uk/embl/>
- European Bioinformatics Institute (EBI), England
- Established in 1980 by the European Molecular Biology Laboratory, Heidelberg, Germany
- Now part of **ENA**, the European Nucleotide Archive, <http://www.ebi.ac.uk/ena/>
- **DDBJ**, <http://www.ddbj.nig.ac.jp/>
- National Institute of Genetics, Japan

Together they form

- International Nucleotide Sequence Database Collaboration, <http://www.insdc.org/>



CENTERFO
R BIOLOGI
CAL SEQU
ENCEANA
LYSIS CBS

FASTA format


```
>alpha-D
ATCTCTGACCGACTCTGCAAGAGAGCTGCTCTGAGGTTGGGGAAGGTGATCCGCCAC
CCAGACTGTGGAGCCGAGGCCCTGGAGAGGTGCGGGCTGAGCTTGGGGAAACCATGGGCA
AGGGGGGGGACTGGGTGGGAGCCCTACAGGGCTGCTGGGGGTGTTTGGCTGGGGGTGAG
CACTGACCATCCCGCTCCCGAGCTGTTCCACACTACCCCGACCAAGACTACTTCC
CCCACTTCGACTTGACCATGGCTCCGACAGGTCCGCAACCAAGCAAGAGGTGTTGG
CCGCTTGGGCAAGCTGTTCAGAGGCTGGGCAACTCAGCCAAAGCCTGTCTGACCTCA
GCGACTGCGATGCTCAACCTGGGTGTGAGCCCTGTCAACTTCAAGGCAAGGCGGGGAC
GGGGGTCAAGGGGCGGGGAGTTGGGGGCGAGGAGCTGTTGGGGATCCGGGGCCATGCC
GGCGGTACTGAGCCCTGTTTGGCTTGGAGCTGCTGGGCGAGTGTCTTCACTGGTGGCTG
GCCACACACTGGGCAAGACTACACCCGGAGGACACTGTGCTTTCGCAAGTTCCTG
TGGGCTGTGTGACCGTGTGGCGAGAGTACAGATAA

>alpha-A
ATGTTCTGTCTTCCCAAGCAAGAGCACTGAAAGCCCTCTTCGCAAAATCGGCGGC
CAGGCCGAGTACTGGGTGTGAGAGCCCTGGAGAGGTATGTGTCTATCCGTATTACCCC
ATCTCTGTGTGTGTGTGATCCATCCATCTGCCCCACTACTTCCCGTCCATAACTG
TCCCTGTCTATGTGGCCCTGGCTGTGTCTGTCTGTCCCAACTGTCCCTGATTGCTC
TGTCCCCAGGTGTTTATCACTACCCCGACCAAGACTACTTCCCGACTTGAGCC
TGTACATGGCTCCGCTCAGATCAAGGGGCAAGAGGTGGCGAGGCACTGGTTG
AGGCTGCCAAGCAATGATGACATCGCTGGTGTCTTCCAGTCAAGGAGGAGTCCACG
CCCAAGGCTCGTGTGGAGCCCGCTCACTTCAAAGTGAGCTTGGGAGAGGGGTGACCA
GTCTGGCTCCCGCTTGGACACACCTCTGGTACCCCTCACTCAACCCCTTGTCTACC
ATCTCTCTTTGGCTTTCAGTGTGGGTCACTGTCTCTGGTGGTGGCGGTCACTT
CCCTCTCTCTGAGCCCGAGGTCCATGCTTCCCTGGCAAGTCTGTGTGGCGTGGG
CACGCTCTTACTGCCAAGTACGTTAA
```

(Handout)

CENTERFO
R BIOLOGI
CAL SEQU
ENCEANA
LYSIS CBS

GenBank format



- Originates from the GenBank database.
- Contains both a DNA sequence and annotations of features (e.g. location of genes).

(handout)

CENTERFO
R BIOLOGI
CAL SEQU
ENCEANA
LYSIS CBS

GenBank format - HEADER

```
LOCUS       CMGLOAD               1185 bp    DNA     linear   VRT 18-APR-2005
DEFINITION  Cairina moschata (duck) gene for alpha-D globin.
ACCESSION   X01831
VERSION     X01831.1  GI:62724
KEYWORDS    alpha-globin; globin.
SOURCE      Cairina moschata (Muscovy duck)
  ORGANISM  Cairina moschata
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Archosauria; Aves; Neognathae; Anseriformes; Anatidae; Cairina.
REFERENCE   1  (bases 1 to 1185)
AUTHORS     Erbil,C. and Niessing,J.
TITLE       The primary structure of the duck alpha D-globin gene: an unusual
            5' splice junction sequence
JOURNAL     EMBO J. 2 (8), 1339-1343 (1983)
PUBMED      10872328
COMMENT     Data kindly reviewed (13-NOV-1985) by J. Niessing.
```

GenBank format - ORIGIN section

```

CENTERFO
RBIOLGDI
CALSFUO
ENCENARA
LYSIS C8S

ORIGIN
1      ctgcgtgggcc  tcagcccccct  caccctccca  cgcctgataag  ataaaggccag  ggcgggaagcg
61      cagggtgcta  taagagtgct  gcccccgggg  ttctctccac  acagaataacc  ctgaattgcc
121     agcctgcacac  gccctgctcg  ccatcgtgac  cgcgcaggac  aagaagctcta  tcgtgcaggt
181     gctgggaagaa  gtggctggcc  accaggagga  attccgaagt  gaagctctgc  agaggtgtgg
241     gctgggaccca  gggggcactc  acagggtggc  cagcaggagga  caggagccct  gcaagcggtg
301     tgggctggact  cccagagcgc  caccgggtgg  ggggtcagagt  ggggaacaaga  cgaggcgacc
361     aaaaactgact  ggctcgctc  cggcaggatg  ttctccgctt  acccccagac  caaacactca
421     ttcceccact  tgaactgcta  cctccggctc  gaacaggctc  gtggccgctc  caagaagact
481     gcggctgccc  tgggaacatg  cgtgaagagc  ttggacaacc  ccagcaggag  cctgctctgag
541     ctccagcaac  ctgcaactgc  ctgcaactgc  cagcactcgt  ggcacactgg  gcaagcggtg
601     gactagaggtc  tctgggcttc  ggggctctag  ggtgtggggt  gcagggtctg  ggtgtccagg
661     ggctctgatt  tccctggatt  tgcactctgt  cggcgagctg  gcccaggctg  ctgtctgtct
721     ggttaccagg  gtccctgggg  cagccagcca  caggcagggg  gctgggattt  catctgggat
781     gctggggcaga  gctcgggatt  gtgtttgaa  tgggaagctg  cgaggggcta  gggccagggt
841     gggggagctca  ggggctcagg  gggagctcgg  gggggtagta  gggagactca  ggggcactat
901     tcogagagag  ggggtactaa  cctcggtttg  ctctcgactc  gctgcagacg  tgcctccagg
961     tgggtgctgc  cgcactgact  ggcacaaagt  acagccggct  gctgctgctg  gctcttgaca
1021    agttctcttc  cgcctgctgc  cgcgtctggt  tccaagaata  cagatgagac  ctctgcctgc
1081    cctctgcacc  ttcataaaga  acaacctaac  cacagctctg  tgtctgtgtg  ttgctgggact
1141    gggcatcggg  ggtcccaggg  agggctgggt  ttgcttcaca  catcc

//

```

GenBank format - FEATURE section

FEATURE	Location/Qualifiers
source	1..1185 /organism="Cairina moschata" /mol_type="genomic DNA" /db_xref="taxon:8855"
CAAT_signal	20..24
TATA_signal	69..73
precursor_mRNA	101..1114 /note="primary transcript"
exon	101..234 /number=1
CDS	join(143..234,387..591,935..1067) /codon_start=1 /product="alpha D-globin" /protein_id="CAAB0866.2" /db_xref="GI:1445876" /db_xref="GDA:PG2003" /db_xref="InterPro:IPR0009971" /db_xref="InterPro:IPR002138" /db_xref="InterPro:IPR002140" /db_xref="InterPro:IPR009050" /db_xref="UniProt/TranS-Prote/P02003" /translation="MLTAEKKKILVQWEKVAQDEEFSEALRQMLFATPQTCTTFF RPLQFQSRQGRDFVVALGDAVSYLMSQAPSESSAHNLVDFDPYFNPFLKA CTCTYLAAGGGVSYFMNAFTFSSDASVAVTAEET"
repeat_region	227..246 /note="direct repeat 1"
intron	235..386 /number=1
repeat_region	289..309 /note="direct repeat 1"
exon	387..591 /number=2
intron	592..939 /number=2
exon	940..1114 /number=3
polyA_signal	1095..1100
polyA_signal	1114

- The exercise guide is linked from the course programme.
- Read the guide carefully - it contains a lot of information about GenBank.
- Remember your handouts:
 - GenBank & FASTA format
 - Eukaryotic gene structure